

CHAP 4: STATISTIQUE  
BAYÉSIENNE

Ces notes de cours sont inspirés de plusieurs supports, notamment des notes de M. Ribatet.

Ref. bibliographie: C.P. Robert - The Bayesian Choice : A Decision-theoretic Motivation - Springer

### I. Introduction

→ On se place dans le cadre d'un modèle statistique paramétrique  $\{X, \mathcal{P}, \{\mathbb{P}_\theta : \theta \in \Theta\}\}$ .  
Le statistique fréquentiste suppose l'existence d'un "vrai" paramètre  $\theta_0$  qui génère les données, i.e.  $(X_1, \dots, X_n) \sim g(\cdot; \theta_0)$  ( $g$  est la densité de probabilité).

Puis on cherche un estimateur  $\hat{\theta}_n$  de  $\theta_0$  qui a de bonnes propriétés (comme EMV):

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \Sigma)$$

L'estimateur  $\hat{\theta}_n$  nous donne une estimation ponctuelle pour  $\theta_0$ , que l'on pourra compléter par une estimation ensembliste (IC) en utilisant ces résultats asymptotiques. Par exemple,  $IC_{\alpha}(0) = [\hat{\theta}_n - q_{1-\frac{\alpha}{2}}^{(0,1)} \hat{\sigma}_n; \hat{\theta}_n + q_{1-\frac{\alpha}{2}}^{(0,1)} \hat{\sigma}_n]$  avec  $\hat{\sigma}_n$  un estimateur de l'écart-type de  $\hat{\theta}_n$ .

En résumé,  $\hat{\theta}_n$  est une v.a. pour laquelle on ne connaît pas que la loi asymptotique

→ Dans le cadre bayésien, on ne considère plus le paramètre  $\theta$  comme déterministe et inconnu - On considère que  $\theta$  est lui-même une v.a.

## II - Fondamentaux de la statistique bayésienne.

En bayésien, on modélise donc également l'incertitude que l'on a sur le paramètre par une loi sur ce paramètre. On se sert de l'expérience accumulée (via les observations  $x_1, \dots, x_n$  recueillies) pour mettre à jour cette incertitude.

### ① - Loi a priori

C'est la loi choisie pour modéliser l'incertitude sur le paramètre. Cette loi a priori permet d'intégrer par exemple un avis d'expert, d'encoder nos connaissances a priori et nos ignorances sur  $\theta$ . AVANT d'observer des données.

On la note en général  $\Pi$ , donc  $\theta \sim \Pi$ .

→ Ex: pour un modèle gaussien:  $X_i \sim \mathcal{N}(\mu, \sigma^2)$ , on pourrait par exemple choisir pour  $\theta = (\mu, \sigma^2)$  la loi  $\Pi(\theta) = \Pi(\mu) \times \Pi(\sigma^2) = \mathcal{N}(\mu_0, T) \times \text{InvGamma}(\alpha, \beta)$

Définition: Les paramètres de la loi a priori (ici  $\mu_0, T, \alpha$  et  $\beta$ ) sont appelés des hyperparamètres. Ces hyperparamètres ont des valeurs fixées par le statisticien, on ne les estime pas!

Définition: Un modèle bayésien est la donnée, pour une v.e. (ou une suite de v.e.) d'une loi conditionnelle et d'une loi a priori:

$$X \sim f(X|\theta)$$

$$\theta \sim \Pi$$

## ② - Loi jointe et loi a posteriori.

(2)

Etant donné que  $\theta$  admet une loi (a priori), on peut exprimer la loi jointe pour un modèle statistique donné, à savoir

$$\Pi(x, \theta) = f(x|\theta) \pi(\theta).$$

Rq: en prévision, la loi jointe n'a que peu d'intérêt en statistique bayésienne. On s'intéresse davantage à la loi a posteriori.

Définition (Loi a posteriori): On peut séparer les cas en 4 catégories:

1) La loi de  $X$  et la loi de  $\theta$  sont discrètes: la loi a posteriori vaut

$$P(\theta = \theta_i | X = x) = \frac{P(X=x | \theta = \theta_i) P(\theta = \theta_i)}{P(X=x)} = \frac{P(X=x | \theta = \theta_i) P(\theta = \theta_i)}{\sum_k P(X=x | \theta = \theta_k) P(\theta = \theta_k)}$$

2) La loi de  $X$  est discrète et la loi de  $\theta$  est continue. Alors on a:

$$\Pi(\theta | X=x) = \frac{P(X=x | \theta) \pi(\theta)}{\int_{\theta \in \Theta} P(X=x | \theta) \pi(\theta) d\theta}$$

3) La loi de  $X$  est continue et la loi de  $\theta$  est discrete. On obtient

$$P(\theta = \theta_i | x) = \frac{f(x | \theta = \theta_i) \pi(\theta = \theta_i)}{\sum_k f(x | \theta = \theta_k) \pi(\theta = \theta_k)}$$

4) La loi de  $X$  est continue et la loi de  $\theta$  aussi. Alors on a:

$$\Pi(\theta | x) = \frac{f(x | \theta) \pi(\theta)}{\int_{\theta \in \Theta} f(x | \theta) \pi(\theta) d\theta}$$

Rq: On résume souvent ces 4 équations par le cas continu/continu; c-e-d:

qu'on appelle loi a posteriori la loi dont la densité est donnée par

$$\pi(\theta|x) = \frac{f(x|\theta) \pi(\theta)}{\int_{\Theta} f(x|\theta) \pi(\theta) d\theta}$$

→ Le dénominateur correspond à la loi marginale de  $x$ . On le note souvent  $m(x)$ . Ce dénominateur joue le rôle de constante de normalisation pour la loi a posteriori : en effet, il est indépendant de  $\theta$  et pour refléter une densité.

→ Pour déterminer le comportement de la loi a posteriori, on travaille souvent à cette constante près, i.e.

$$\pi(\theta|x) \propto \frac{f(x|\theta) \pi(\theta)}{\int_{\Theta} f(x|\theta) \pi(\theta) d\theta}$$

- Rq:
- Si on cherche le maximum de la loi a posteriori, il est facile de calculer la loi marginale.
  - $x$  est le vecteur des valeurs observées en pratique, permettant d'intégrer l'expérience et de mettre à jour la loi "a priori" qui devient "a posteriori".
  - On donne le nom de stéthoscopie bayésienne car le lien avec le Théorème de Bayes est évident. Exercice - Annex Chap 2 ex. 12, 14

### III - Plus loin sur la loi a priori

#### ① - Comment la choisir ?

C'est souvent l'étape oubliée en stéthoscopie bayésienne.

(3)

Pourtant, c'est une étape fondamentale, et c'est même ce qui constitue la plus grande différence avec la statistique fréquentiste.

En général, ce choix peut être motivé par plusieurs types de considération :

- basé sur des expériences similaires du passé,
- basé sur un avis d'expert ou une intuition,
- basé sur la faisabilité des calculs (même si aujourd'hui les méthodes par simulation de type MCMC permettent de contourner cet écueil),
- basé sur la volonté de n'apporter aucune info nouvelle pouvant biaiser l'estimation.

## ② - Loi a priori conjuguée

Définition : Une famille  $\mathcal{F}$  de lois de probabilité sur  $\Theta$  est dite conjuguée pour le modèle statistique  $\{P(x|\theta) : x \in X, \theta \in \Theta\}$  si, pour tout  $\Pi \in \mathcal{F}$ , la loi a posteriori  $\Pi(\theta|x) \propto f(x|\theta) \Pi(\theta)$  appartient aussi à  $\mathcal{F}$ .

→ Exemple: modèle Poisson-Gamma: 
$$\begin{aligned} & X_i \sim P(\theta) \\ & \Rightarrow \text{Faire le calcul...} \quad \theta \sim \text{Gamma}(\alpha, \lambda) \end{aligned}$$

→ Intérêt principal: choisir une loi a priori conjuguée permet d'obtenir une loi a posteriori explicite. En pratique, la mise à jour se fait à travers les paramètres de la loi a priori, ce qui facilite aussi l'interprétation du changement dû à l'accumulation de l'expérience.

→ Remarque: Une loi conjuguée peut être déterminée en considérant la forme de la vraisemblance  $f(x|\theta)$  et en prenant une loi a priori de la même forme (vue comme une fonction du paramètre).

→ Quelques exemples de lois conjuguées:

$f(x \theta)$	$\pi(\theta)$	$\pi(\theta x)$
$\mathcal{N}(\theta, \sigma^2)$	$\mathcal{N}(\mu, \tau^2)$	$\mathcal{N}\left(\frac{\sigma^2\mu + \tau^2 x}{\sigma^2 + \tau^2}, \frac{\sigma^2\tau^2}{\sigma^2 + \tau^2}\right)$
$\mathcal{P}(\theta)$	$\text{Gamma}(x, \lambda)$	$\text{Gamma}(\lambda + x, \lambda + 1)$
$\text{Gamma}(r, \theta)$	$\text{Gamma}(x, \lambda)$	$\text{Gamma}(x + r, \lambda + x)$
$\mathcal{B}(n, \theta)$	$\text{Beta}(x, \beta)$	$\text{Beta}(x + \alpha, \beta + n - x)$
$\mathcal{U}(\mu, \frac{1}{\theta})$	$\text{Gamma}(x, \lambda)$	$\text{Gamma}\left(\lambda + \frac{1}{2}, \lambda + \frac{(\mu - x)^2}{2}\right)$
Pareto	Gamma	

→ L'impact de la loi a priori s'amerise quand  $n \rightarrow \infty$ . (Illustration ?)

### ③. Loi a priori non-informative.

C'est le cas où l'on a peu (ou pas) d'information sur  $\theta$ . On cherche donc à ce que les données parlent d'elles-mêmes au maximum, et de ne pas influencer la loi a posteriori. On étudie ici 2 familles de lois a priori non-informatives: la loi de Laplace et celle de Jeffreys. (le cas le + évident étant la loi uniforme...)

#### a) Loi de Laplace:

Definition : la loi a priori de Laplace consiste à poser  $\pi(\theta) \propto \prod_{\theta \in \Theta}$

C'est donc la loi uniforme (discrete ou continue selon ④), voire la mesure de Lebesgue dans le cas d'une loi impropre.

→ Cette loi présente toutefois quelques défauts → elle peut mener à une loi a priori impropre → elle est non-invariante par reparamétrisation.

→ Ex: prouver le modèle Exponentiel ( $\lambda$ ),  $\lambda > 0$ ; ainsi que sa reparamétrisation

$$\lambda = \exp(\theta), \theta \in \mathbb{R}.$$

• Dans le premier cas, on a  $\pi_1(\lambda) \propto \prod_{\lambda > 0}$

• Dans le 2<sup>e</sup> cas, on a  $\pi_2(\theta) \propto 1 \Rightarrow \pi_2(\lambda) \propto \frac{1}{\lambda} \pi_2(\ln(\lambda)) \prod_{\lambda > 0} \propto \frac{1}{\lambda} \prod_{\lambda > 0}$

⇒ 2<sup>e</sup> cas fortement informatif: la notion de non-informat dépend du paramètre  $\lambda$

## b) A priori de Jeffreys:

(4)

On a vu qu'une bonne notion de loi a priori non-informatrice est une loi invariante par reparamétrisation. (ce qui n'est pas forcément le cas d'une loi Uniforme en fonction de la paramétrisation du problème).

L'a priori de Jeffreys est fondé sur l'information de Fisher.

Définition : La loi a priori de Jeffreys est donnée par

$$\pi(\theta) \propto |I(\theta)|^{1/2}, \text{ où } |A| \text{ est le déterminant de } A.$$

- Cette loi est invariante par reparamétrisation.
- Elle peut conduire à des lois a priori impropre.
- Elle est à éviter quand  $\dim(\mathcal{A}) > 1$ .
- $I(\theta)$  est un indicateur de la quantité d'information apportée par le modèle  $f(x|\theta)$ : il est donc intuitif que les valeurs de  $\theta$  pour lesquelles  $I(\theta)$  est plus grande doivent être plus probables a priori.
- Remarque sur la reparamétrisation: soit  $\phi = h(\theta)$  avec  $h$  un  $C^1$ -diffeomorphisme.  
En notant  $\tilde{\pi}$  la loi a priori de  $\theta$ , alors  $\phi$  est de loi  $\tilde{\pi}$  avec  
 $\tilde{\pi}(\phi) = \pi(h^{-1}(\phi)) |(h^{-1})'(\phi)|$  - De plus on a  $\tilde{I}(\phi) = I(\phi) |(h^{-1})'(\phi)|^2$ , donc  
 $\tilde{\pi}(\phi) \propto \sqrt{\tilde{I}(\phi)}$  (cela justifie la racine carré).
- Ex?

## IV - Inférence bayésienne

Le résultat du travail bayésien donne accès à la loi a posteriori  $\Pi(\theta|x)$ .  
Cette loi contient beaucoup d'information sur  $\theta$  (beaucoup plus qu'en fréquentiste), et en pratique on résume cette loi à des indicateurs plus simples, comme :

- le maximum a posteriori (MAP),
- la médiane a posteriori,
- la moyenne a posteriori,
- un quantile a posteriori d'ordre  $p$ .

### ① - Estimateur bayésien

Définition : Si  $d(\cdot, \cdot)$  est une fonction de perte, le risque bayésien de l'estimateur  $\hat{\theta} = t(x)$  du paramètre  $\theta$  pour la perte  $d$  est

$$R^B(\hat{\theta}) = \mathbb{E}_{(\theta, x)}[d(\hat{\theta}, \theta)] = \mathbb{E}_{(\theta, x)}[d(t(x), \theta)]$$

Rq: On peut généraliser cette définition à une fonction de  $\theta$ ,  $\psi(\theta)$ , avec un estimateur de cette fonction de  $\theta$ .

N.B.: On intègre sous la loi jointe  $(\theta, x)$ .

Proposition : Le risque bayésien est la moyenne du risque fréquentiste/fréquentiel suivant la loi a priori  $\Pi$ . Ainsi,

$$R^B(\hat{\theta}) = \mathbb{E}_{\theta \sim \Pi}[R(\hat{\theta}, \theta)] = \int R(\hat{\theta}, \theta) \Pi(\theta) d\theta$$

Nous allons maintenant aborder le sujet d'un estimateur optimal au sens du risque bayésien.

Définition : Étant donné une fonction de perte  $d(\cdot, \cdot)$ , l'estimation bayésienne

associée à cette fonction de perte est définie, pour tout  $x \in X$ ,

$$\text{par } t^B(x) = \underset{\theta}{\operatorname{argmin}} E[d(t, \theta) | X=x].$$

L'estimateur bayésien est donc  $\hat{\theta}^B = t^B(x)$ .

Théorème : L'estimateur bayésien est l'estimateur qui minimise le risque bayésien.

Proposition : L'estimateur bayésien associé au risque quadratique est donné par l'espérance a posteriori:  $\hat{\theta}^B = t^B(x) = E_{\pi}[\Theta | X]$ .

Rq: On retrouve le fait que pour minimiser un critère de moindre carré, la meilleure estimation est la moyenne (cf régression linéaire).

- Si le risque est le risque  $L^1$  (perte en valeur absolue), l'estimateur bayésien est donné par la médiane de la loi a posteriori  $\pi(\cdot | X)$ .
- On peut généraliser tous ces résultats à une fonction de  $\theta$ , en remplaçant  $\Theta$  par  $\varphi(\theta)$  partout. Exercices - Exercice Chap 3. ex 14, 15, 17, 18, 20

## ② - Intervalle de crédibilité.

La notion est proche de celle de l'intervalle de confiance, mais elle est ≠.

En statistique fréquentiste, on rappelle qu'un intervalle de confiance a des bornes aléatoires. Le niveau  $(1-\alpha)\%$  de cet intervalle correspond à la proportion  $(1-\alpha)\%$  que  $n$  réalisations de cet IC contiennent le vrai paramètre lorsque  $n \rightarrow \infty$ .

## a) Régrion $\alpha$ -crédible:

Définition: Pour une loi a priori  $\Pi$  donnée, un ensemble  $G_x \subset \mathbb{R}$  est un ensemble  $\alpha$ -crédible si

$$P_{\Pi}(\theta \in G_x | x) \geq 1-\alpha.$$

→ Interprétation:  $P_{\Pi}(\theta \in G_x | x) = \int_{\mathbb{R}} \mathbb{1}_{\{\theta \in G_x\}} \underbrace{\Pi(\theta | x)}_{\text{loi a posteriori}} d\theta$

On intègre donc l'expérience observée  $x$  dans la définition de cet intervalle - loi a posteriori.

## b) Intervalle de crédibilité

On simplifie les choses en considérant que  $\theta$  est scalaire.

Définition: Pour une loi a priori  $\Pi$  donnée, un intervalle  $IC_x \subset \mathbb{R}$  est un intervalle de crédibilité de niveau  $1-\alpha$  si

$$P_{\Pi}(\theta \in IC_x | x) = 1-\alpha.$$

→ Souvent on utilise des intervalles de crédibilité symétriques, c'est :

$$IC_x = \left[ q_{\Pi}\left(\frac{\alpha}{2}, x\right), q_{\Pi}\left(1-\frac{\alpha}{2}, x\right) \right], \text{ avec}$$

$$q_{\Pi}(p, x) = \inf \{u \in \mathbb{R} : P_{\Pi}(\theta < u | x) \geq p\}.$$

→ Comparaison entre intervalle de confiance / intervalle de crédibilité:

Exercice, Annexe Chap. 2

ex 13

•  $IC$  est aléatoire

• La prochaine réalisation de  $I$  sera avec  $(1-\alpha)\%$  de chance de contenir  $\theta_0$ .

•  $\theta$  est aléatoire

• Ayant observé  $x$ , il y a  $(1-\alpha)\%$  de chance que  $IC_x$  contienne  $\theta_0$ .

### ③ - Loi prédictive a posteriori:

On dispose de  $x = (x_1, \dots, x_n)$ .

On souhaite prédire la prochaine observation  $x_{n+1}$ .

→ En fréquentiste, on a tendance à utiliser le prédicteur  $E[X]$ , avec  $X \sim f(\cdot; \theta)$  et  $\hat{\theta}$  estimateur de  $\theta$ .

⇒ On ne tient donc pas compte de l'incertitude d'estimation sur  $\theta$ .

Définition : On appelle loi prédictive a posteriori la loi de densité

$$\Pi(x_{n+1} | x) = \int f(x_{n+1} | \theta) \Pi(\theta | x) d\theta.$$

Comme en fréquentiste, on utilise ensuite l'espérance sur cette loi pour définir le prédicteur bayésien :

$$\hat{x}_{n+1} = \int x_{n+1} \Pi(x_{n+1} | x) d_{x_{n+1}}$$

## II - Simulations

Enfin, pour conclure ce chapitre, il est essentiel de mentionner que des méthodes génératrices performantes (mais coûteuses en calcul) existent pour approcher la distribution a posteriori du paramètre  $\theta$  (la loi  $\Theta | X$ ) dans un cadre général (pas forcément une loi a priori conjuguée, ...).

C'est des techniques dépassant le cadre de ce cours, mais l'étudiant intéressé pourra notamment approfondir :

- les techniques basées sur la simulation Monte Carlo
- les algorithmes plus spécifiques de type MCMC :
  - Algorithme de Metropolis-Hastings,
  - Algorithme de Gibbs (échantillonneur)
- les modèles bayésiens hiérarchiques.